

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)**ScienceDirect**

Procedia Computer Science 78 (2016) 276 – 283

**Procedia**  
Computer Science

International Conference on Information Security & Privacy (ICISP2015), 11-12 December 2015,  
Nagpur, INDIA

## A Novel Algorithmic Approach for an Automatic Data Placement for NUMA Based Design

Jitendra Madarkar, V.Gopi Chand, P.V.Sai Kiran Reddy, Harsh Arora

*SCSE VIT University vellore, Tamilnadu, 632014, India*

### Abstract:

Scalability is a key concern for SMP based architecture in the current context. NUMA based architecture design seems to be a promising hope addressing this key concern. At the same time CC-NUMA based design architecture demands a deeper understanding and open vistas for key areas of improvement. Our proposed research tries to investigate, evolve and analyze one of the key design issues for NUMA machine and proposes an innovative solution to address this key design issue under investigation in the current phase of our work.

© 2016 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Peer-review under responsibility of organizing committee of the ICISP2015

**Keywords:** NUMA; SM; CC-NUMA.

### 1. INTRODUCTION

With the need for multiprocessing design, shared memory architecture provides processors to share a common memory. Uniform memory access (UMA) is one such architecture where the memory is uniformly accessed by all the processors. If the data size is increased, the processing speed need to be increased where in turn number of processors should be more. Since UMA shares a common bus, the allocation of bandwidth to the processors is a bottleneck which results in a scalability problem.

To overcome this problem, Non-Uniform Memory Access (NUMA) based design has been proposed. In NUMA, the memory is physically distributed among NUMA nodes where a common global address space is maintained<sup>12</sup>. These nodes are connected with an interconnect. Memory access is fast and there is no issue of bandwidth. In UMA the latencies are uniform whereas in NUMA it differs with the distance of NUMA nodes. In NUMA each node is associated with a local memory<sup>13</sup>. If the memory accesses happen from this local node, then there is no issue of bandwidth. As there can be an access to global address space at some time, there is a need to access data from the remote nodes. This creates a latency difference between local and remote memory accesses and hence remote access should be minimized. Even though NUMA overcomes the scalability issue, some factors need to be optimized. Such key design factors include Data placement, Processor Affinity, Load Balancing, Cache Coherence, Thread Scheduling<sup>14</sup>. In this paper we plan to discuss different strategies to ensure optimal data placement in NUMA based system using automatic page migration and replication scheme.

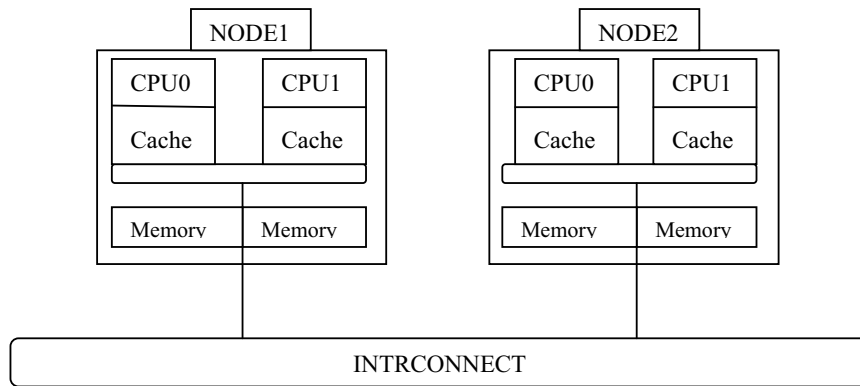


Fig.1.: NUMA (Non Uniform Memory Architecture)

In NUMA, design for effective data placement is the major concern. Improper data locality leads to increase in remote access. For an efficient Data placement proper memory management along with process scheduling has to be done. Key requirements for design of optimal data placement NUMA engine require:

- Placer data shall have to be available locally to a NUMA node.
- Placer should avoid the overload transmission through interconnect.
- Remote access has to be minimized.
- If a particular node has insufficient memory, memory reclaim mechanism needs to be built in the Data placer.
- In case of memory reclaim failure, there should be an alternate mechanism with placer to allocate neighbouring NUMA nodes (preferred nodes) memory.

Memory affinity is one such data placement policy wherein we can achieve the above mentioned factors. Memory affinity involves data allocation, data migration, data replication which is efficiently being designed/proposed in this paper.

The rest of the paper is planned as follows. In section 1, we provided the various ways of providing memory affinity in NUMA architecture, section 3 gives a detailed explanation of various memory policies followed up with our algorithm in section 4. Section 5 summarizes our design along with discussion and section 6 provides the Conclusion with relevant future work in future.

## 2. LITERATURE SURVEY

To assure memory affinity in NUMA machines several works have been adopted. These works lead to some major solutions. In this section we briefly discuss various methodologies to assure memory affinity. The cons and pros of each method was identified and discussed.

### 2.1 Using LibNUMA

LibNUMA supports page migration, memory policies, and CPU bindings using kernel system calls through. These system calls allow programmer to allocate memory in run time. It majorly includes "numa\_migrate\_pages()", "numa\_set\_mpolicy()", "numa\_move\_pages()", "mbind()" and "numa\_get\_mpolicy()" [3]. The main advantage of this kernel system call is that memory allocation can be controlled in a better way [4].

Usage of LibNUMA is a complex task since it involves bit masks, pointers, memory pages. Also to use this LibNUMA, one must manually enter the system calls (numactl) so that the application customs to the memory policies, which becomes a complicated work for the programmers. So there must be a mechanism by which the application must automatically adapt to the architecture. One such proposed mechanism is discussed in the later sections.

### 2.2 Open MP

Open MP is a language which supports efficient parallelism. The application must be coded using OpenMP

primitives in order to get multithreading. To place the data in better way in NUMA machines, certain extensions were added to the OpenMP and it requires an explicit support from compilers. But all compilers don't give the provision for OpenMP and directives are applied in a static manner<sup>9</sup>.

In<sup>5</sup> a mechanism to guarantee memory affinity on NUMA machine using OpenMP was presented by the author. The main idea in the paper was to make relation between threads and data. Also the work provides some suggestions of how to extend OpenMP in NUMA machines. Their results proved that OpenMP can perform well on tightly-coupled NUMA machines. Their work was not extended to automatic data placement. Our proposed work tries to provide this mechanism using automatic memory affinity in run time.

In<sup>6</sup> author presented an efficient memory allocation in OpenMP using certain set of OpenMP directives. These directives guide developers to allocate data efficiently on the NUMA architecture. All these directives are limited to the FORTRAN language. Using these directives efficient data distribution and page locality can be achieved<sup>7,8</sup>. But these directives have to be included for an application in an explicit manner which requires hardware specifications in prior. There must be a better approach independent of prior knowledge on hardware which must provide efficient data placement.

### 3. BASIS WORK

#### *Kernel's Memory affinity policies*

Memory affinity is ensuring that processing unit always has their data close to them. So, to achieve utmost performance on NUMA machine the number of remote access during the execution has to minimize. Processor and data need to schedule so as the distance between them are to be closed. To minimize the distance and accordingly increase the performance of NUMA architecture we need some mechanism and tool in order to solve memory allocation, replication and migration to ensure memory affinity in NUMA system. However, none of these solutions provide portability as memory affinity control. To achieve this, memory affinity is ensured by applying a memory policy for an entire process. First touch algorithm was proposed in Linux 2.6.24 kernel which was the default policy to manage memory affinity. It places the page with respect to thread on the NUMA node that access it first. By this policy data is allocated by thread or master thread to its local memory which reduces remote accesses. However first touch policy apply on application that access symmetric data (regular). If thread does not access similar data, it leads to high number of remote access.

To reduce remote access the page should be migrated between NUMA nodes to make sure that data is closer to the processor that access it. Before page migration the following sequence of step has to be followed

- a. The victim page which is to be migrated should be removed from the Least Recently used lists.
- b. All the references of PTE (page table entries) to the old page are freed and the page is put into sleep till the migration of the page is over.
- c. If suppose any kernel references are made to the page, then the page migration is aborted and should be further processed after the kernel usage is over.
- d. New kernel references are to be made for the new page after migration.

Next touch policy was introduced to provide automatic page migration. This policy is applicable for regular and irregular data. Next touch policy provide automatic dynamic page migration using "copy on write" (COW). Using COW page migration occurs only it is essentially needed.

Generally page table entry (PTE) contain protection bit (write access and read access bit) to check the page detail. COW is implemented in LINUX to modify the write access bit of the page from the page table entry which indirectly generates a page fault on a write access.

The figure 2 explains the flow of process occurred in the next touch policy. Initially by using first touch policy, thread associates/binds page with respective NUMA node. In next touch policy the page is marked by flag bit and indicates that it will be used in near future. In order to mark the page using flag bit in LINUX it is accomplished using **madvise()** system call.

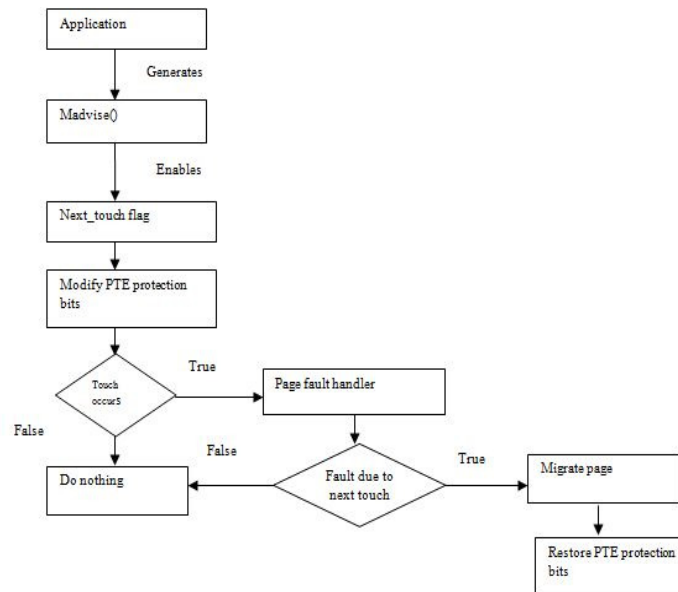


Fig.2.: Next touch policy implementation

Whenever we want to set the next touch flag, we have to modify the PTE protection bit. Due to the load balancing there is need to migrate the threads from one NUMA node to another irrespective of the data location of the thread. This leads to page fault whenever touch occurs for the page. When this happen page fault handler take care of the faults by checking the next touch flag of the page. It also sees that actual page fault is occurred because of Next Touch policy or a real page Fault .If the fault is due to the Next touch, then page will be migrated. After the page migration, the Next touch flag is removed from the corresponding buffer and enables its original protection bits in the **PTE**. Same procedure happens if Next touch fault happens.

In general some pages are protected. This can be done by using **mprotect ()** system call which prevents application to access that page. If that page is to be accessed, it leads to segmentation fault which is handled by custom signal handler<sup>4</sup>. This custom handler removes the protection temporarily, then migrates the resultant buffers and restores its protection back. This mechanism is not flexible to implement as we are calling the **mprotect()**(system call) twice to handle the segmentation fault, resulting in changing the TLBs frequently.

The next touch policy produces a better performance when compared with **numa\_migrate\_pages()**(libNUMA's API) because **numa\_migrate\_pages()** Shifts entire process address space on to new NUMA node which is not necessary as discussed in the previous sections.

The proposed design addressed all the above mentioned policies along with replication in next touch policy.

#### 4. Design and architecture for proposed system

Our proposed system is trying to provide a replication mechanism for the aforementioned next touch policy with some threshold constraint. Our algorithm tries to give an effective mechanism for Memory affinity which targets in reducing the remote accesses. The proposed mechanism also tries to incorporate all the LibNUMA support (migration, replication, preferred and interleaved) *automatically*

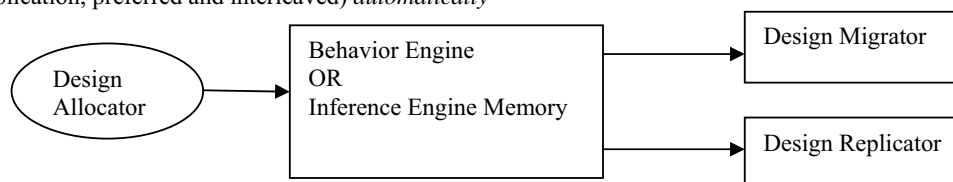


Fig.3.: Design and Architecture of the proposed system

**Module 1: Design Allocator**

Design allocator allocates memory to a node by default policy (First touch). According to first touch policy, memory will be allocated to the thread that first touches it. It sets the behavior of the shared page. Then it sets the next touch flag indicating that the page will be accessed in the future. Before doing this we must check whether the memory is available for migrating or replicating the pages to the destination nodes or not. If page fault value is less than the threshold value and memory is available, migrate the page. In case if the memory is not available then memory reclaim has to be done by using **madvise\_free()**

**Module 2: Behavior engine**

Behavior engine decides whether page candidate is fit to Migrate or Replicate based on the threshold value. The threshold value is to be checked with page fault counter. It also provides a access mechanism for pages (e.g. More write access implicitly disables the page from replication)

**Module 3: Design Migrator or Replicator**

This forms the core enabler for Migration or Replicator. Before the migration or replication happens again we have to check for sufficient memory availability. If the sufficient memory is not present memory reclaim (**madvise\_free()**) has to be happened. If the page fault count is greater than threshold frequency data replication has to be done. Page fault counter should be incremented only on the next touch. Whereas if the page fault count is less than threshold frequency data migration has to be done. Figure 4 depicts the overall proposed methodology flow discussed above.

To include replication mechanism we introduced some variables as threshold, reset (RST), and write frequency threshold. The threshold frequency checks the no of access of the page for a certain period across then NUMA nodes. The reset flag resets the page fault counter value after a particular clock cycle time. The write frequency threshold finds how many write access happened to a page. If it is high, then check the write frequency threshold, if the number of write access is greater than write frequency threshold, do nothing. Else check for the memory again and replicate the page.

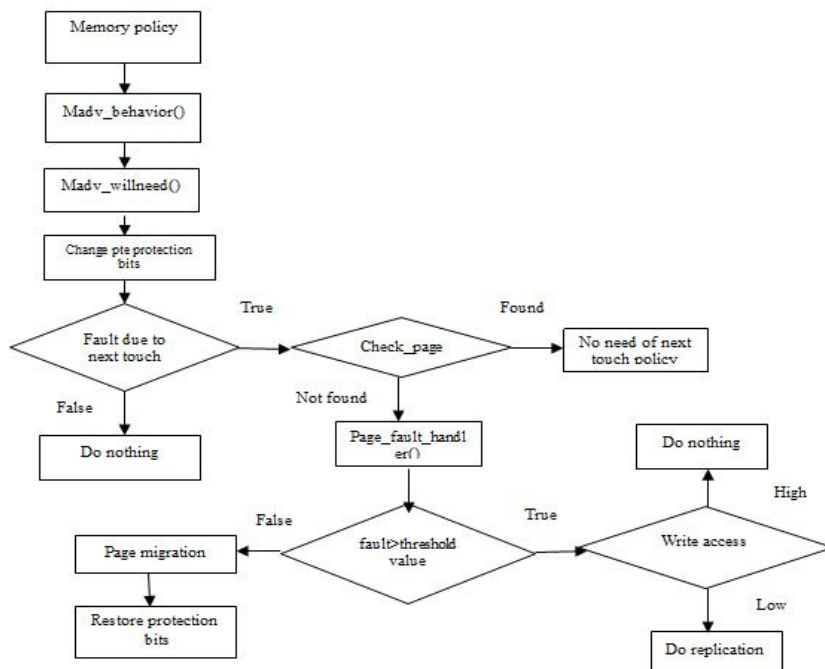


Fig. 4.: Next touch policy with replication

**Algorithm:****Init\_configuration\_cum\_allocator\_NUMA()****Step 1:** Allocate a page in NUMA node as follows:

- Step1.1.a Choose default **memory policy()**
- Step 1.1.b Initialize the Memory Access parameters (like **page\_table\_entry**, **protectionbits**, **pid**, **bitmask**, **writeaccess** **bit**, **vma**, **main memory** etc...)
- Step1.2 set the behavior of the shared space using **madv\_behavior()**
- Step 1.3 set the next touch flag using **madv\_willneed()**
- Step1.4 Change the page table entry protection bits i.e **pte\_modify(page table entry, get\_protection(0))**.
- Step1.5 Initialize threshold and reset value

**Shared behavior process****Step 2** Check the recent access to Page i.eif (**touch==true**)do the following steps

```

{
    /*checking whether the page fault occurred from the next touch or not*/
    2.1.Check the pages for its local or remote access using
    Page_found=check_page(pid,bitmask,pageaddress)
    2.2. If page is being from local node i.e if (Page_found==1)
        Do nothing and abort ()
        /* next touch policy not needed*/
    Else
        /* page fault handler mechanism and increment the page fault count */
    2.2. Generate a page fault using
    intpage_fault_count=page_fault_handler(addr,writeaccess,pageableentry ,vma ,mm)
        2.2.1Increment the page fault using page_fault_count++;
        /*until reset values doesn't cross the periodical time value, if crosses reset the page fault values*/
}

```

**Step 3 Replication or Migration**

3.1. Check the Page faults limit&amp; compare the same with threshold limit set using.

If (**page\_fault\_count >= threshold value**)

```

{
    3.2.Replicate the page using page_replicate(node mask_all,vma,page table entry,pd,mm)
}
else
/* (if page_fault_count < threshold value)
{
    3.1 Migrate the page using next touch policy as page_migrate(mm,pte,ptl,mm,vma)
    /*provide migration*/
}

```

**Algorithm for Main design block design for data placer()****Automatic\_data\_placer\_for\_NUMA()**

```

{
    Step a) Configure the NUMA nodes &allocate the pages;
    Init_configuration_cum_allocator_NUMA ( );/* step 1 above*/
    Step b)Organize the NUMA shared process behavior;
        Perform step2 above;
    Step c) Based upon the threshold frequency value
    Either Migrate or Replicate; /* step 3* above */
}

```

The algorithm describes the next touch policy with replication. Initially every application obeys the default memory policy as first touch. **Madvise()** system call gives a prior knowledge to kernel of how an application expects to use the memory. The following table describes the different functionalities of **madvise()** system call. By using **madv\_willneed()** system call for which we change the page table entry bits(read/write access flag), indicates that page will be used in the future. Now whenever fault occurs, kernel confirms that this fault is by next

<i>System call</i>	<i>Description</i>
<b><i>madvise_normal</i></b>	<i>Default kernel way of accessing the addresses.</i>
<b><i>madvise_sequential</i></b>	<i>Informs kernel that application will access a listed range addresses in a successive way.</i>
<b><i>madvise_random</i></b>	<i>Tells kernel that page references are in a random manner.</i>
<b><i>madvise_willneed</i></b>	<i>The specified address range will be referenced in future.</i>
<b><i>madvise_dontneed</i></b>	<i>The specified address range will not be referenced in future.</i>
<b><i>madvise_free</i></b>	<i>Intimating kernel that the specified range of addresses are no longer important(these addresses are freed when the memory pressure is high)</i>

touch policy and decides to take decision of migration or replication.

Table 1: Different functionalities of madvise()

## 5. Comparative analysis and Discussions

In this proposed design we have investigated and addressed the key concerns or limitations imposed by the existing works. Our Algorithmic design proposed seems to be outperforming with respect to LibNUMA specifically.

Consider an application which is generating 200 threads, where these threads are randomly distributed in parallel among the NUMA nodes. In LibNUMA using **mbind()** we can allocate memory to any node. **cpubind()** allows process to execute on a set of CPUs for a specified node. LibNUMA also provides page migration using **numa\_migrate\_pages()** where a large buffers are to be moved which leads to increase in latencies (unnecessary pages are also transferred) [15]. For the LibNUMA support the programmer should alter the application code manually which is a complex task [9].

By first touch policy adopted by us, the thread which touches the page first will be allocated to the corresponding node to make the access local. If other threads from different nodes need to access the same page (remote) very frequently in a short period, then migration becomes a tedious job where in TLB values are to be flushed frequently. So instead of migration we replicate the page based on the page fault count (if page fault is greater than the threshold count) in our proposed system addressing the key concern.

Replication may also cause memory congestion. Duplicate copies indeed result in redundancy of data. So in order to make the memory available we use **madvise(madv\_free)** resets in a periodical interval and the replicated page is removed based on the necessity. For every migration or replication the status of the page is checked in advance. Hence we ensured that our proposed system provides a better alternative to the existing works with automatic processing built-in.



## 6. Conclusion and future work

Scalability being a major concern in UMA based system; NUMA addressed this key concern very effectively. With the advent of many core environments where memory is distributed among the different cores, it is challenging to design a thread scheduler along with proper data distribution across different nodes in an efficient manner.

We have thoroughly investigated existing solution with respect to data placement and explored their limitation in current context (lack of dynamic adaption to run time directives in OpenMP and lack of automatic placer in the current LibNUMA APIs).

We have proposed, investigated and evolved an efficient design for NUMA based machine for automatic data placement in scalable fashion, there by addressing the key concern (dynamic migration and replication)

In next phase of our work we plan to validate our proposed framework for NUMA data placement engine and performance analyze the same with respect to existing aforementioned works. Subsequently we plan to investigate the other key areas of improvement identified erstwhile in NUMA based design.

## References

1. I. C. Compiler, . Thread affinity interface,. 2010.  
<http://software.intel.com/en-us/intel-compilers/>
2. G. C. Compiler, . Thread affinity interface,. 2010.  
<http://gcc.gnu.org/onlinedocs/libgomp/Environment-Variables.html>
3. Andi Kleen SUSE Labs . An NUMA API for Linux. .
4. A. Kleen, . A NUMA API for Linux,. Tech.Rep. Novell-4621437, April 2005.  
Online. . Available:<http://whitepapers.zdnet.co.uk/0,1000000651,260150330p,00.htm>
5. D. S. Nikolopoulos, E. Artiaga, E. Ayguadé, and J. Labarta, . Exploiting Memory Affinity in OpenMP Through Schedule Reuse,. SIGARCH Computer Architecture News, vol. 29, no. 5, pp. 49–55, 2001.
6. J. Bircsak, P. Craig, R. Crowell, Z. Cvetanovic, J. Harris, C. A. Nelson, and C. D. Offner, . Extending OpenMP for NUMA Machines,. in SC '00: Proceedings of the 2000 ACM/IEEE Conference on Supercomputing, Dallas, Texas, USA, 2000.
7. H. Richardson, . High Performance Fortran: history, overview and current developments,. Tech. Rep. TMC-261, 1996.  
<http://hpff.rice.edu/publications/index.html>
8. S. Benkner and T. Brandes, . Efficient Parallel Programming on Scale Shared Memory Systems with High Performance Fortran,. Concurrency: Practice and Experience, vol. 14, pp. 789–803, 2002.
9. Christiane Pousa Ribeiro, Jean-François Méhaut . MAi: Memory Affinity Interface. inria-00344189, version 6 - 14 Jun 2010
10. Charles Koelbel, David Loveman, Robert Schreiber, Guy Steele, and Mary Zosel. *The High Performance Fortran Handbook*, 1994.
11. François Broquedis, Nathalie Furmento, Brice Goglin, Raymond Namyst, Pierre-André Wacrenier . Dynamic Task and Data Placement over NUMA Architectures: an OpenMP Runtime Perspective. published in . International Workshop
12. T. Mu, J. Tao, M. Schulz, and S. A. McKee, . Interactive Locality Optimization on NUMA Architectures,. in SoftVis '03: Proceedings of the 2003 ACM Symposium on Software Visualization. New York, NY, USA: ACM, 2003, pp. 133–ff.
13. J. Marathe and F. Mueller, . Hardware Profile-Guided Automatic Page Placement for ccNUMA Systems,. in PPOPP '06: Proceedings of the eleventh ACM SIGPLAN symposium on Principles and practice of parallel programming. New York, NY, USA: ACM, 2006, pp. 90–99.  
Online. . Available: <http://portal.acm.org/citation.cfm?id=1122987>
14. Christoph Lameter. Local and Remote Memory: Memory in a Linux/NUMA System. In Linux Symposium (OLS2006), Ottawa, Canada, July 2006.
15. Christian Terboven, Dieter an Mey, Dirk Schmidl, Henry Jin, and Thomas Reichstein. Data and Thread Affinity in OpenMP Programs. In Proceedings of the 2008 workshop on Memory access on future processors (MAW '08), pages 377–384, New York, NY, 2008. ACM.